**(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)**

**(72) Inventors; and**
**(75) Inventors/Applicants** *(for US only)*: **PENN, Sharron, G.** [GB/US]; 617 South Delaware Street, San Mateo, CA 94402 (US). **HANZEL, David, K.** [US/US]; 968 Loma Verde Avenue, Palo Alto, CA 94303 (US). **CHEN, Wen-sheng** [CN/US]; 210 Easy Street #25, Mountain View, CA 94043 (US). **RANK, David, R.** [US/US]; 117 El Dorado Commons, Fremont, CA 94539 (US).

**(74) Agent: RONNING, Royal, N., Jr.**; Amersham Pharmacia Biotech, Inc., 800 Centennial Avenue, P.O. Box 1327, Piscataway, NJ 08855 (US).

**(54) Title:** HUMAN GENOME-DERIVED SINGLE EXON NUCLEIC ACID PROBES USEFUL FOR ANALYSIS OF GENE EXPRESSION IN HUMAN HEART

**(57) Abstract:** A single exon nucleic acid microarray comprising a plurality of single exon nucleic acid probes for measuring gene expression in a sample derived from human heart is described. Also described are single exon nucleic acid probes expressed in the heart and their use in methods for detecting gene expression.

HUMAN GENOME-DERIVED SINGLE EXON NUCLEIC ACID PROBES USEFUL
FOR ANALYSIS OF GENE EXPRESSION IN HUMAN HEART

CROSS REFERENCE TO RELATED APPLICATIONS

5

The present application is a continuation-in-part of U.S.
patent application serial nos. 09/632,366, filed August 3,
2000 and 09/608,408, filed June 30, 2000; claims the
benefit under 35 U.S.C. s 119(e) of U.S. provisional patent
10   application serial nos. 60/236,359, filed September 27,
2000, 60/234,687, filed September 21, 2000, 60/207,456,
filed May 26, 2000, and 60/180,312, filed February 4, 2000;
and further claims the benefit under 35 U.S.C. s 119(a) of
UK patent application no. 0024263.6, filed October 4, 2000,
15   the disclosures of which are incorporated herein by
reference in their entireties.

REFERENCE TO SEQUENCE LISTING AND INCORPORATION BY
REFERENCE THEREOF

20

The present application includes a Sequence Listing in
electronic format, filed pursuant to PCT Administrative
Instructions 801 - 806 on a single CD-R disc, in
triplicate, containing a file named pto_HEART.txt, created
25   24 January 2001, having 20,186,946 bytes.  The Sequence
Listing contained in said file on said disc is incorporated
herein by reference in its entirety.

Field of the Invention

30

        The present invention relates to genome-derived
single exon microarrays useful for verifying the expression
of regions of genomic DNA predicted to encode protein. In
particular, the present invention relates to unique genome-
35   derived single exon nucleic acid probes expressed in human

heart and single exon nucleic acid microarrays that include
such probes.


Background of the Invention

5        For almost two decades following the invention of
general techniques for nucleic acid sequencing, Sanger *et
al.*, *Proc. Natl. Acad. Sci. USA* 70(4):1209-13 (1973);
Gilbert *et al.*, *Proc. Natl. Acad. Sci. USA* 70(12):3581-4
(1973), these techniques were used principally as tools to
10    further the understanding of proteins — known or
suspected — about which a basic foundation of biological
knowledge had already been built.  In many cases, the
cloning effort that preceded sequence identification had
been both informed and directed by that antecedent
15    biological understanding.
         For example, the cloning of the T cell receptor
for antigen was predicated upon its known or suspected cell
type-specific expression, by its suspected membrane
association, and by the predicted assembly of its gene via
20    T cell-specific somatic recombination. Subsequent
sequencing efforts at once confirmed and extended
understanding of this family of proteins.  Hedrick *et al.*,
*Nature* 308(5955):153-8 (1984).
         More recently, however, the development of high
25    throughput sequencing methods and devices, in concert with
large public and private undertakings to sequence the human
and other genomes, has altered this investigational
paradigm: today, sequence information often precedes
understanding of the basic biology of the encoded protein
30    product.
         One of the approaches to large-scale sequencing
is predicated upon the proposition that expressed
sequences — that is, those accessible, through isolation of
mRNA — are of greatest initial interest.  This "expressed
35    sequence tag" ("EST") approach has already yielded vast

amounts of sequence data (see for example Adams *et al.*, *Science* 252:1651 (1991); Williamson, *Drug Discov. Today* 4:115 (1999)). For nucleic acids sequenced by this approach, often the only biological information that is

5    known *a priori* with any certainty is the likelihood of biologic expression itself. By virtue of the species and tissue from which the mRNA had originally been obtained, most such sequences are also annotated with the identity of the species and at least one tissue in which expression

10   appears likely.

More recently, the pace of genomic sequencing has accelerated dramatically. When genomic DNA serves as the initial substrate for sequencing efforts, expression cannot be presumed; often the only *a priori* biological information

15   about the sequence includes the species and chromosome (and perhaps chromosomal map location) of origin.

With the ever-accelerating pace of sequence accumulation by directed, EST, and genomic sequencing approaches — and in particular, with the accumulation of

20   sequence information from multiple genera, from multiple species within genera, and from multiple individuals within a species — there is an increasing need for methods that rapidly and effectively permit the functions of nucleic sequences to be elucidated. And as such functional

25   information accumulates, there is a further need for methods of storing such functional information in meaningful and useful relationship to the sequence itself; that is, there is an increasing need for means and apparatus for annotating raw sequence data with known or

30   predicted functional information.

Although the increase in the pace of genomic sequencing is due in large part to technological changes in sequencing strategies and instrumentation, Service, *Science* 280:995 (1998); Pennisi, *Science* 283: 1822-1823 (1999),

35   there is an important functional motivation as well.

While it was understood that the EST approach would rarely be able to yield sequence information about the noncoding portions of the genome, it now also appears the EST approach is capable of capturing only a fraction of

5    a genome's actual expression complexity.

For example, when the *C. elegans* genome was fully sequenced, gene prediction algorithms identified over 19,000 potential genes, of which only 7,000 had been found by EST sequencing. *C. elegans* Sequencing Consortium,

10   *Science* 282:2012 (1998). Analogously, the recently completed sequence of chromosome 2 of *Arabidopsis* predicts over 4000 genes, Lin *et al.*, *Nature*, 402:761 (1999), of which only about 6% had previously been identified via EST sequencing. Although the human genome has the

15   greatest depth of EST coverage, it is still woefully short of surrendering all of its genes. One recent estimate suggests that the human genome contains more than 146,000 genes, which would at this point leave greater than half of the genes undiscovered. It is now predicted that many

20   genes, perhaps 20 to 50%, will only be found by genomic sequencing.

There is, therefore, a need for methods that permit the functional regions of genomic sequence — and most importantly, but not exclusively, regions that

25   function to encode genes — to be identified.

Much of the coding sequence of the human genome is not homologous to known genes, making detection of open reading frames ("ORFs") and predictions of gene function difficult. Computational methods exist for predicting

30   coding regions in eukaryotic genomes. Gene prediction programs such as GRAIL and GRAIL II, Uberbacher *et al.*, *Proc. Natl. Acad. Sci. USA* 88(24):11261-5 (1991); Xu *et al.*, *Genet. Eng.* 16:241-53 (1994); Uberbacher *et al.*, *Methods Enzymol.* 266:259-81 (1996); GENEFINDER, Solovyev *et*

35   *al.*, *Nucl. Acids. Res.* 22:5156-63 (1994); Solovyev *et al.*,

4

*Ismb* 5:294-302 (1997); and GENESCAN, Burge *et al.*, *J. Mol. Biol.* 268:78-94 (1997), predict many putative genes without known homology or function. Such programs are known, however, to give high false positive rates. Burset *et al.*,

5   *Genomics* 34:353-367 (1996). Using a consensus obtained by a plurality of such programs is known to increase the reliability of calling exons from genomic sequence. Ansari-Lari *et al.*, *Genome Res.* 8(1):29-40 (1998)

Identification of functional genes from genomic

10  data remains, however, an imperfect art. For example, in reporting the full sequence of human chromosome 21, the Chromosome 21 Mapping and Sequencing Consortium reports that prior bioinformatic estimates of human gene number may need to be revised substantially downwards. *Nature*

15  405:311-199 (2000); Reeves, *Nature* 405:283-284 (2000).

Thus, there is a need for methods and apparatus that permit the functions of the regions identified bioinformatically — and specifically, that permit the expression of regions predicted to encode protein — readily

20  to be confirmed experimentally.

Recently, the development of nucleic acid microarrays has made possible the automated and highly parallel measurement of gene expression. *Reviewed in* Schena (ed.), DNA Microarrays : A Practical Approach

25  (Practical Approach Series), Oxford University Press (1999) (ISBN: 0199637768); *Nature Genet.* 21(1)(suppl):1 - 60 (1999); Schena (ed.), Microarray Biochip: Tools and Technology, Eaton Publishing Company/BioTechniques Books Division (2000) (ISBN: 1881299376).

30  It is common for microarrays to be derived from cDNA/EST libraries, either from those previously described in the literature, such as those from the I.M.A.G.E. consortium, Lennon *et al.*, *Genomics* 33(1):151-2 (1996), or from the construction of "problem specific" libraries

35  targeted at a particular biological question, R.S. Thomas

*et al.*, *Cancer Res.* (in press).  Such microarrays by definition can measure expression only of those genes found in EST libraries, and thus have not been useful as probes for genes discovered solely by genomic sequencing.

5　　　　　　The utility of using whole genome nucleic acid microarrays to answer certain biological questions has been demonstrated for the yeast *Saccharomyces cerevisiae*.  De Risi *et al.*, *Science* 278:680 (1997).  The vast majority of yeast nuclear genes, approximately 95% however, are single

10　exon genes, *i.e.*, lack introns, Lopez *et al.*, *RNA* 5:1135-1137 (1999); Goffeau *et al.*, *Science* 274:563-67 (1996), permitting coding regions more readily to be identified. Whole genome nucleic acid microarrays have not generally been used to probe gene expression from more complex

15　eukaryotic genomes, and in particular from those averaging more than one intron per gene.

　　　　　　Diseases of the heart and vascular system are a significant cause of human morbidity and mortality. Increasingly, genetic factors are being found that

20　contribute to predisposition, onset, and/or aggressiveness of most, if not all, of these diseases.  Although mutations in single genes have on occasion been identified as causative, these disorders are for the most part believed to have polygenic etiologies. There is a need for methods

25　and apparatus that permit prediction, diagnosis and prognosis of diseases of the human heart, particularly those diseases with polygenic etiology.


Summary of the Invention

30

　　　　　　The present invention solves these and other problems in the art by providing methods and apparatus for predicting, confirming, and displaying functional information derived from genomic sequence.  The present

35　invention also provides apparatus for verifying the

expression of putative genes identified within genomic
sequence.

In particular, the invention provides novel
genome-derived single exon nucleic acid microarrays useful
5   for verifying the expression of putative genes identified
within genomic sequence.

The present invention also provides compositions
and kits for the ready production of nucleic acids
identical in sequence to, or substantially identical in
10  sequence to, probes on the genome-derived single exon
microarrays of the present invention.

Accordingly, in a first aspect of the invention,
there is provided a spatially-addressable set of single
exon nucleic acid probes for measuring gene expression in a
15  sample derived from human heart, comprising a plurality of
single exon nucleic acid probes according to any one of the
nucleotide sequences set out in SEQ ID NOs: 1 - 9,980 or a
complementary sequence, or a portion of such a sequence.

By plurality is meant at least two, suitably at
20  least 20, most suitably at least 100, preferably at least
1000 and, most preferably, upto 5000.

In one embodiment of the first aspect, each of
said plurality of probes is separately and addressably
amplifiable.

25         In an alternative embodiment, each of said
plurality of probes is separately and addressably
isolatable from said plurality.

In a preferred embodiment, each of said plurality
of probes is amplifiable using at least one common primer.
30  Preferably, each of said plurality of probes is amplifiable
using a first and a second common primer.

In yet another embodiment, said set of single
exon nucleic acid probes comprises between 50 - 20,000
probes, for example, 50 - 5000.

35         Suitably, said set of single exon nucleic acid

7

probes comprises at least 50 - 1000 discrete single exon nucleic acid probes having a sequence as set out in any of SEQ ID NOS.: 1 - 19,771 or a complimentary sequence, or a portion of such a sequence.

Preferably, the average length of the single exon nucleic acid probes is between 200 and 500 bp. It is preferred that the average length should be at least 200bp, suitably at least 250bp, most suitably at least 300bp, preferably at least 400bp and, most preferably, 500 bp.

In another embodiment, the single exon nucleic acid probes lack prokaryotic and bacteriophage vector sequence. It is preferred that at least 50%, suitably at least 60%, most suitably at least 70%, preferably at least 75%, more preferably at least 80, 85, 90, 95 or 99% of said single exon nucleic acid probes lack prokaryotic and bacteriophage vector sequence.

In another preferred embodiment, said single exon nucleic acid lack homopolymeric stretches of A or T. It is preferred that at least 50%, suitably at least 60%, most suitably at least 70%, preferably at least 75%, more preferably at least 80, 85, 90, 95 or 99% of said single exon nucleic acid probes lack homopolymeric stretches of A or T.

Preferably, a spatially-addressable set of single exon nucleic acid probes in accordance with the first aspect of the invention is is addressably disposed upon a substrate.

Suitable substrates include a filter membrane which may, preferably, be nitrocellulose or nylon. The nylon may preferably, be positively-charged. Other suitable substrates include glass, amorphous silicon, crystalline silicon, and plastic. Further suitable materials include polymethylacrylic, polyethylene, polypropylene, polyacrylate, polymethylmethacrylate, polyvinylchloride, polytetrafluoroethylene, polystyrene, polycarbonate,

8

polyacetal, polysulfone, celluloseacetate,
cellulosenitrate, nitrocellulose, and mixtures thereof.

      In a second aspect of the invention, there is
provided a microarray comprising a spatially addressable
5   set of single exon nucleic acid probes in accordance with
the first aspect of the invention.

      In one embodiment, a genome-derived single-exon
microarray is packaged together with such an ordered set of
amplifiable probes corresponding to the probes, or one or
10  more subsets of probes, thereon.  In alternative
embodiments, the ordered set of amplifiable probes is
packaged separately from the genome-derived single exon
microarray.

      In another aspect, the invention provides genome-
15  derived single exon nucleic acid probes useful for gene
expression analysis, and particularly for gene expression
analysis by microarray. In particular embodiments of this
aspect, the present invention provides human single-exon
probes that include specifically-hybridizable fragments of
20  SEQ ID Nos. 9,981 – 19,771, wherein the fragment hybridizes
at high stringency to an expressed human gene.  In
particular embodiments, the invention provides single exon
probes comprising SEQ ID Nos. 1 – 9,980.

      Accordingly, in a third aspect of the invention,
25  there is provided a single exon nucleic acid probe for
measuring human gene expression in a sample derived from
human heart which is a nucleic acid molecule comprising a
nucleotide sequence as set out in any of SEQ ID NOs.: 1 –
9,980 or a complementary sequence or a fragment thereof
30  wherein said probe hybridizes at high stringency to a
nucleic acid expressed in the human heart.

      In one embodiment, a single exon nucleic acid
probe in accordance with the third aspect comprises a
nucleotide sequence as set out in any of SEQ ID NOs.: 9,981
35  – 19,771 or a complementary sequence or a fragment thereof.

In a fourth aspect of the invention, there is provided a single exon nucleic acid probe for measuring human gene expression in a sample derived from human heart which is a nucleic acid molecule having a sequence encoding a peptide comprising a peptide sequence as set out in any of SEQ ID NOs.: 19,772 - 29,119 or a complementary sequence or a fragment thereof wherein said probe hybridizes at high stringency to a nucleic acid expressed in the human heart.

Preferably, a single exon nucleic acid probe in accordance with the third or fourth aspects of the invention comprises between at least 15 and 50 contiguous nucleotides of said SEQ ID NO:. It is preferred that the single exon nucleic acid probe comprises at least 15, suitably at least 20, more suitably at least 25 or preferably at least 50 contiguous nucleotides of said SEQ ID NO:.

In another preferred embodiment, a single exon nucleic acid probe in accordance with the third or fourth aspects of the invention is between 3kb and 25kb in length. It is preferred that said probe is no more than 3kb, suitably no more than 5kb, more suitably no more than 10kb, preferably 15kb, more preferably 20kb or, most preferably, no more than 20kb in length.

Preferably, a single exon nucleic acid probe in accordance with either the fifth or sixth aspect of the invention is DNA, preferably single-stranded DNA, RNA or PNA.

In another embodiment of either the third or fourth aspect of the invention, a single exon nucleic acid probe is detectably labeled. Suitable detectable labels include a radionuclide, a fluorescent label or a first member of a specific binding pair. Suitable fluorescent labels include dyes such as cyanine dyes, preferably Cy3 and Cy5 although other suitable dyes will be known to those skilled in the art.

In a particularly preferred embodiment, a single exon nucleic acid probe in accordance with either the third or fourth aspect of the invention lacks prokaryotic and bacteriophage vector sequence. In yet another embodiment, a

5 single exon nucleic acid probe in accordance with either the third or fourth aspect of the invention lacks homopolymeric stretches of A or T.

In a fifth aspect of the invention, there is provided an amplifiable nucleic acid composition,

10 comprising:

the single exon nucleic acid probe in accordance with either of the third or fourth aspects of the invention; and at least one nucleic acid primer;

wherein said at least one primer is sufficient to

15 prime enzymatic amplification of said probe.

In an sixth aspect of the invention, there is provided a method of measuring gene expression in a sample derived from human heart, comprising:

contacting the single exon microarray in

20 accordance with the second aspect of the invention, with a first collection of detectably labeled nucleic acids, said first collection of nucleic acids derived from mRNA of human heart; and then

measuring the label detectably bound to each

25 probe of said microarray.

In a seventh aspect of the invention, there is provided a method of identifying exons in a eukaryotic genome, comprising:

algorithmically predicting at least one exon from

30 genomic sequence of said eukaryote; and then

detecting specific hybridization of detectably labeled nucleic acids to a single exon probe,

wherein said detectably labeled nucleic acids are derived from mRNA from the heart of said eukaryote, said

35 probe is a single exon probe having a fragment identical in

sequence to, or complementary in sequence to, said
predicted exon, said probe is included within a single exon
microarray in accordance with the first aspect of the
invention, and said fragment is selectively hybridizable at
5    high stringency.

In a eighth aspect of the invention, there is
provided a method of assigning exons to a single gene,
comprising:

identifying a plurality of exons from genomic
10   sequence in accordance with the seventh aspect of the
invention; and then

measuring the expression of each of said exons in
a plurality of tissues and/or cell types using
hybridization to single exon microarrays having a probe
15   with said exon,

wherein a common pattern of expression of said
exons in said plurality of tissues and/or cell types
indicates that the exons should be assigned to a single
gene.

20   In an ninth aspect of the invention, there is
provided a nucleic acid sequence as set out in any of SEQ
ID NOs: 1 - 19,771 wherein said sequence encodes a peptide.

In a tenth aspect of the invention, there is
provided a peptide encoded by a sequence comprising a
25   sequence as set out in any of SEQ ID NOs: 9,981 - 19,771,
or a complementary sequence or coding portion thereof.

In a preferred embodiment, a peptide may be
encoded by a sequence comprising a sequence set out in any
of SEQ ID NOS.: 1 - 9,980.

30   In a further aspect, the invention provides
peptides comprising an amino acid sequence translated from
the DNA fragments, said amino acid sequences comprising SEQ
ID NOS.: 9,981 - 19,771.

Accordingly in a eleventh aspect of the invention
35   there is provided a peptide comprising a sequence as set

12

out in any of SEQ ID NOs: 19,772 - 29,119, or fragment thereof.

In another aspect, the invention provides means for displaying annotated sequence, and in particular, for displaying sequence annotated according to the methods and apparatus of the present invention. Further, such display can be used as a preferred graphical user interface for electronic search, query, and analysis of such annotated sequence.

## Detailed Description of the Invention

### Definitions

As used herein, the term "microarray" and phrase "nucleic acid microarray" refer to a substrate-bound collection of plural nucleic acids, hybridization to each of the plurality of bound nucleic acids being separately detectable. The substrate can be solid or porous, planar or non-planar, unitary or distributed.

As so defined, the term "microarray" and phrase "nucleic acid microarray" include all the devices so called in Schena (ed.), DNA Microarrays: A Practical Approach (Practical Approach Series), Oxford University Press (1999) (ISBN: 0199637768); Nature Genet. 21(1)(suppl):1 - 60 (1999); and Schena (ed.), Microarray Biochip: Tools and Technology, Eaton Publishing Company/BioTechniques Books Division (2000) (ISBN: 1881299376). As so defined, the term "microarray" and phrase "nucleic acid microarray" further include substrate-bound collections of plural nucleic acids in which the nucleic acids are distributably disposed on a plurality of beads, rather than on a unitary planar substrate, as is described, inter alia, in Brenner et al., Proc. Natl. Acad. Sci. USA 97(4):166501670 (2000); in such case, the term "microarray" and phrase "nucleic

acid microarray" refer to the plurality of beads in aggregate.

As used herein with respect to a nucleic acid microarray, the term "probe" refers to the nucleic acid that is, or is intended to be, bound to the substrate; in such context, the term "target" thus refers to nucleic acid intended to be bound thereto by Watson-Crick complementarity. As used herein with respect to solution phase hybridization, the term "probe" refers to the nucleic acid of known sequence that is detectably labeled.

As used herein, the expression "probe comprising SEQ ID NO.", and variants thereof, intends a nucleic acid probe, at least a portion of which probe has either (i) the sequence directly as given in the referenced SEQ ID NO., or (ii) a sequence complementary to the sequence as given in the referenced SEQ ID NO., the choice as between sequence directly as given and complement thereof dictated by the requirement that the probe hybridize to mRNA.

As used herein, the term "open reading frame" and the equivalent acronym "ORF" refer to that portion of an exon that can be translated in its entirety into a sequence of contiguous amino acids i.e. a nucleic acid sequence that, in at least one reading frame, does not possess stop codons; the term does not require that the ORF encode the entirety of a natural protein.

As used herein, the term "amplicon" refers to a PCR product amplified from human genomic DNA, containing the predicted exon.

As used herein the term "exon" refers to the consensus prediction of the various exon and gene predicting algorithms i.e. a nucleic acid sequence bioinformatically predicted to encode a portion of a natural protein.

As used herein, the term "peptide" refers to a sequence of amino acids. The sequences referred to as

14

PEPTIDE SEQ ID NOS.: are the predicted peptide sequences that would be translated from one of the exons, or a portion thereof set out in exon SEQ ID NOS.:. The codons encoding the peptide are wholly contained within the exon.

5      As used herein, a "portions" of a defined nucleotide sequence or sequences can be and, preferably, are fragments unique to that sequence or to one or a combination of those sequences. A fragment unique to a nucleic acid molecule is one that is a signature for the

10     larger nucleic acid molecule.

As used herein, the phrase "expression of a probe" and its linguistic variants means that the ORF present within the probe, or its complement, is present within a target mRNA.

15     As used herein, "stringent conditions" refers to parameters well known to those skilled in the art. When a nucleic acid molecule is said to be hybridisable to another of a given sequence under "stringent conditions" it is meant that it is homologous to the given sequence.

20     As used herein, the phrase "specific binding pair" intends a pair of molecules that bind to one another with high specificity. Binding pairs are said to exhibit specific binding when they exhibit avidity of at least $10^7$, preferably at least $10^8$, more preferably at least $10^9$

25     liters/mole. Nonlimiting examples of specific binding pairs are: antibody and antigen; biotin and avidin; and biotin and streptavidin.

As used herein with respect to the visual display of annotated genomic sequence, the term "rectangle" means

30     any geometric shape that has at least a first and a second border, wherein the first and second borders each are capable of mapping uniquely to a point of another visual object of the display.

As used herein, a "Mondrian" means a visual

35     display in which a single genomic sequence is annotated

15

with predicted and experimentally confirmed functional
information.


5    Brief Description of the Drawings


         The present invention is further illustrated with
reference to the following non-limiting figures and
examples in which:
10           FIG. 1 illustrates a process for predicting
functional regions from genomic sequence, confirming the
functional activity of such regions experimentally, and
associating and displaying the data so obtained in
meaningful and useful relationship to the original sequence
15   data;
         FIG. 2 further elaborates that portion of the
process schematized in FIG. 1 for predicting functional
regions from genomic sequence;
         FIG. 3 illustrates a Mondrian visual display;
20           FIG. 4 presents a Mondrian showing a hypothetical
annotated genomic sequence;
         FIG. 5 is a histogram showing the distribution of
ORF length and PCR products as obtained, with ORF length
shown in black and PCR product length shown in dotted
25   lines;
         FIG. 6 is a histogram showing the distribution,
among exons predicted according to the methods described,
of expression as measured using simultaneous two color
hybridization to a genome-derived single exon microarray.
30   The graph shows the number of sequence-verified products
that were either not expressed ("0"), expressed in one or
more but not all tested tissues ("1" - "9"), or expressed
in all tissues tested ("10");
         FIG. 7 is a pictorial representation of the
35   expression of verified sequences that showed expression

16

with signal intensity greater than 3 in at least one
tissue, with: FIG. 7A showing the expression as measured by
microarray hybridization in each of the 10 measured
tissues, and the expression as measured "bioinformatically"
5      by query of EST, NR and SwissProt databases; with FIG. 7B
showing the legend for display of physical expression
(ratio) in FIG. 7A; and with FIG. 7C showing the legend for
scoring EST hits as depicted in FIG. 7A;

FIG. 8 shows a comparison of normalized CY3
10     signal intensity for arrayed sequences that were identical
to sequences in existing EST, NR and SwissProt databases or
that were dissimilar (unknown), where black denotes the
signal intensity for all sequence-verified products with a
BLAST Expect ("E") value of greater than 1e-30 (1 x $10^{-30}$)
15     ("unknown") and a dotted line denotes sequence-verified
spots with a BLAST expect ("E") value of less than 1e-30 (1
x $10^{-30}$) ("known");

FIG. 9 presents a Mondrian of BAC AC008172 (bases
25,000 to 130,000), containing the carbamyl phosphate
20     synthetase gene (AF154830.1); and

FIG. 10 is a Mondrian of BAC A049839.


Methods and Apparatus for Predicting, Confirming,
25     Annotating, and Displaying Functional Regions From Genomic
Sequence Data


FIG. 1 is a flow chart illustrating in broad
outline a process for predicting functional regions from
30     genomic sequence, confirming and characterizing the
functional activity of such regions experimentally, and
then associating and displaying the information so obtained
in meaningful and useful relationship to the original
sequence data.

35     The initial input into process 10 of the present

invention is drawn from one or more databases 100
containing genomic sequence data.  Because genomic sequence
is usually obtained from subgenomic fragments, the sequence
data typically will be stored in a series of records
5    corresponding to these subgenomic sequenced fragments.
Some fragments will have been catenated to form larger
contiguous sequences ("contigs"); others will not.  A
finite percentage of sequence data in the database will
typically be erroneous, consisting *inter alia* of vector
10   sequence, sequence created from aberrant cloning events,
sequence of artificial polylinkers, and sequence that was
erroneously read.

          Each sequence record in database 100 will
minimally contain as annotation a unique sequence
15   identifier (accession number), and will typically be
annotated further to identify the date of accession,
species of origin, and depositor.  Because database 100 can
contain nongenomic sequence, each sequence will typically
be annotated further to permit query for genomic sequence.
20   Chromosomal origin, optionally with map location, can also
be present.  Data can be, and over time increasingly will
be, further annotated with additional information, in part
through use of the present invention, as described below.
Annotation can be present within the data records, in
25   information external to database 100 and linked to the
records thereto, or through a combination of the two.

          Databases useful as genomic sequence database 100
in the present invention include GenBank, and particularly
include several divisions thereof, including the
30   htgs(draft), NT (nucleotide, command line), and NR
(nonredundant) divisions.  GenBank is produced by the
National Institutes of Health and is maintained by the
National Center for Biotechnology Information (NCBI).
Databases of genomic sequence from species other than
35   human, such as mouse, rat, Arabidopsis, *C. elegans*, *C.*

18

*brigsii, Drosophila,* zebra fish, and other higher eukaryotic organisms will also prove useful as genomic sequence database 100.

5          Genomic sequence obtained by query of genomic sequence database 100 is then input into one or more processes 200 for identification of regions therein that are predicted to have a biological function as specified by the user. Such functions include, but are not limited to, encoding protein, regulating transcription, regulating

10 message transport after transcription into mRNA, regulating message splicing after transcription into mRNA, of regulating message degradation after transcription into mRNA, and the like. Other functions include directing somatic recombination events, contributing to chromosomal

15 stability or movement, contributing to allelic exclusion or X chromosome inactivation, and the like.

          The particular genomic sequence to be input into process 200 will depend upon the function for which relevant sequence is to be identified as well as upon the

20 approach chosen for such identification. Process step 200 can be iterated to identify different functions within a given genomic region. In such case, the input often will be different for the several iterations.

          Sequences predicted to have the requisite

25 function by process 200 are then input into process 300, where a subset of the input sequences suitable for experimental confirmation is identified. Experimental confirmation can involve physical and/or bioinformatic assay. Where the subsequent experimental assay is

30 bioinformatic, rather than physical, there are fewer constraints on the sequences that can be tested, and in this latter case therefore process 300 can output the entirety of the input sequence.

          The subset of sequences output from process 300

35 is then used in process 400 for experimental verification

and characterization of the function predicted in process 200, which experimental verification can, and often will, include both physical and bioinformatic assay.

Process 500 annotates the sequence data with the
5    functional information obtained in the physical and/or bioinformatic assays of process 400. Such annotation can be done using any technique that usefully relates the functional information to the sequence, as, for example, by incorporating the functional data into the sequence data
10   record itself, by linking records in a hierarchical or relational database, by linking to external databases, by a combination thereof, or by other means well known within the database arts. The data can even be submitted for incorporation into databases maintained by others, such as
15   GenBank, which is maintained by NCBI.

As further noted in FIG. 1, additional annotation can be input into process 500 from external sources 600.

The annotated data is then displayed in process 800, either before, concomitantly with, or after optional
20   storage 700 on nontransient media, such as magnetic disk, optical disc, magnetooptical disk, flash memory, or the like.

FIG. 1 shows that the experimental data output from process 400 can be used in each preceding step of
25   process 10: e.g., facilitating identification of functional sequences in process 200, facilitating identification of an experimentally suitable subset thereof in process 300, and facilitating creation of physical and/or informational substrates for, and performance of subsequent assay, of
30   functional sequences in process 400.

Information from each step can be passed directly to the succeeding process, or stored in permanent or interim form prior to passage to the succeeding process. Often, data will be stored after each, or at least a
35   plurality, of such process steps. Any or all process steps

can be automated.

FIG. 2 further elaborates the prediction of functional sequence within genomic sequence according to process 200.

5       Genomic sequence database 100 is first queried 20 for genomic sequence.

The sequence required to be returned by query 20 will depend, in the first instance, upon the function to be identified.

10      For example, genomic sequences that function to encode protein can be identified *inter alia* using gene prediction approaches, comparative sequence analysis approaches, or combinations of the two.  In gene prediction analysis, sequence from one genome is input into process
15  200 where at least one, preferably a plurality, of algorithmic methods are applied to identify putative coding regions.  In comparative sequence analysis, by contrast, corresponding, *e.g.*, syntenic, sequence from a plurality of sources, typically a plurality of species, is input into
20  process 200, where at least one, possibly a plurality, of algorithmic methods are applied to compare the sequences and identify regions of least variability.

The exact content of query 20 will also depend upon the database queried.  For example, if the database
25  contains both genomic and nongenomic sequence, perhaps derived from multiple species, and the function to be determined is protein coding regions in human genomic sequence, the query will accordingly require that the sequence returned be genomic and derived from humans.

30      Query 20 can also incorporate criteria that compel return of sequence that meets operative requirements of the subsequent analytical method.  Alternatively, or in addition, such operative criteria can be enforced in subsequent preprocess step 24.

35      For example, if the function sought to be

21

identified is protein coding, query 20 can incorporate
criteria that return from genomic sequence database 100
only those sequences present within contigs sufficiently
long as to have obviated substantial fragmentation of any
5   given exon among a plurality of separate sequence
fragments.

Such criteria can, for example, consist of a
required minimal individual genomic sequence fragment
length, such as 10 kb, more typically 20 kb, 30 kb, 40kb,
10  and preferably 50 kb or more, as well as an optional
further or alternative requirement that sequence from any
given clone, such as a bacterial artificial chromosome
("BAC"), be presented in no more than a finite maximal
number of fragments, such as no more than 20 separate
15  pieces, more typically no more than 15 fragments, even more
typically no more than about 10 - 12 fragments.

Results using the present invention have shown
that genomic sequence from bacterial artificial chromosomes
(BACs) is sufficient for gene prediction analysis according
20  to the present invention if the sequence is at least 50 kb
in length, and if additionally the sequence from any given
BAC is presented in fewer than 15, and preferably fewer
than 10, fragments.  Accordingly, query 20 can incorporate
a requirement that data accessioned from BAC sequencing be
25  in fewer than 15, preferably fewer than 10, fragments.

An additional criterion that can be incorporated
into the query can be the date, or range of dates, of
sequence accession.  Although the process has been
described above as if genomic sequence database 100 were
30  static, it is of course understood that the genomic
sequence databases need not be static, and indeed are
typically updated on a frequent, even hourly, basis.  Thus,
as further described in Examples 1 and 2, *infra*, it is
possible to query the database for newly added sequence,
35  either newly added after an absolute date, or newly added

22

relative to a prior analysis performed using the methods
and apparatus of the present invention. In this way, the
process herein described can incorporate a dynamic,
temporal component.

5          One utility of such temporal limitation is to
identify, from newly accessioned genomic sequence, the
presence of novel genes, particularly those not previously
identified by EST sequencing (or other sequencing efforts
that are similarly based upon gene expression). As further

10   described in Example 1, such an approach has shown that
newly accessioned human genomic sequence, when analyzed for
sequences that function to encode protein, readily
identifies genes that are novel over those in existing EST
and other expression databases. This makes the methods of

15   the present invention extremely powerful gene discovery
tools. And as would be appreciated, such gene discovery
can be performed using genomic sequence from species other
than human.

          If query 20 incorporates multiple criteria, such

20   as above-described, the multiple criteria can be performed
as a series of separate queries or as a single query,
depending in part upon the query language, the complexity
of the query, and other considerations well known in the
database arts.

25          If query 20 returns no genomic sequence meeting
the query criteria, the negative result can be reported by
process 22, and process 200 (and indeed, entire process 10)
ended 23, as shown. Alternatively, or in addition to
report and termination of the initial inquiry, a new query

30   20 can be generated that takes into account the initial
negative result.

          When query 20 returns sequence meeting the query
criteria, the returned sequence is then passed to optional
preprocessing 24, suitable and specific for the desired

35   analytical approach and the particular analytical methods

23

thereof to be used in process 25.

Preprocessing 24 can include processes suitable
for many approaches and methods thereof, as well as
processes specifically suited for the intended subsequent
5    analysis.

Preprocessing 24 suitable for most approaches and
methods will include elimination of sequence irrelevant to,
or that would interfere with, the subsequent analysis.
Such sequence includes repetitive sequence, such as Alu
10   repeats and LINE elements, vector sequence, artificial
sequence, such as artificial polylinkers, and the like.
Such removal can readily be performed by identification and
subsequent masking of the undesired sequence.

Identification can be effected by comparing the
15   genomic sequence returned by query 20 with public or
private databases containing known repetitive sequence,
vector sequence, artificial sequence, and other artifactual
sequence.  Such comparison can readily be done using
programs well known in the art, such as CROSS_MATCH, or by
20   proprietary sequence comparison programs the engineering of
which is well within the skill in the art.

Alternatively, or in addition, undesirable,
including artifactual, sequence can be identified
algorithmically without comparison to external databases
25   and thereafter removed.  For example, synthetic polylinker
sequence can be identified by an algorithm that identifies
a significantly higher than average density of known
restriction sites.  As another example, vector sequence can
be identified by algorithms that identify nucleotide or
30   codon usage at variance with that of the bulk of the
genomic sequence.

Once identified, undesired sequence can be
removed.  Removal can usefully be done by masking the
undesired sequence as, for example, by converting the
35   specific nucleotide references to one that is unrecognized

24

by the subsequent bioinformatic algorithms, such as "X".
Alternatively, but at present less preferred, the undesired
sequence can be excised from the returned genomic sequence,
leaving gaps.

5          Preprocessing 24 can further include selection
from among duplicative sequences of that one sequence of
highest quality.  Higher quality can be measured as a lower
percentage of, fewest number of, or least densely clustered
occurrence of ambiguous nucleotides, defined as those
10 nucleotides that are identified in the genomic sequence
using symbols indicating ambiguity.  Higher quality can
also or alternatively be valued by presence in the longest
contig.

          Preprocessing 24 can, and often will, also
15 include formatting of the data as specifically appropriate
for passage to the analytical algorithms of process 25.
Such formatting can and typically will include, *inter alia*,
addition of a unique sequence identifier, either derived
from the original accession number in genomic sequence
20 database 100, or newly applied, and can further include
additional annotation.  Formatting can include conversion
from one to another sequence listing standard, such as
conversion to or from FASTA or the like, depending upon the
input expected by the subsequent process.

25          Preprocessing, which can be optional depending
upon the function desired to be identified and the
informational requirements of the methods for effecting
such identification, is followed by sequence processing 25,
where sequences with the desired function are identified
30 within the genomic sequence.

          As mentioned above, such functions can include,
but are not limited to, encoding protein, regulating
transcription, regulating message transport after
transcription into mRNA, regulating message splicing after
35 transcription, of regulating message degradation, and the

                              25

like. Other functions include directing somatic
recombination events, contributing to chromosomal stability
or movement, contributing to allelic exclusion or X
chromosome inactivation, or the like.

5          The methods of the present invention are
particularly useful for gene discovery, that is, for
identifying, from genomic sequence, regions that function
to encode genes, and in a particularly useful embodiment,
for identifying regions that function to encode genes not

10    hitherto identified by expression-based or directed cloning
and sequencing. In conjunction with verification using the
novel single exon microarrays of the present invention, as
further described below, the methods herein described
become powerful gene discovery tools.

15          Accordingly, in a preferred embodiment of the
present invention, process 25 is used to identify putative
coding regions. Two preferred approaches in process 25 for
identifying sequence that encodes putative genes are gene
prediction and comparative sequence analysis.

20          Gene prediction can be performed using any of a
number of algorithmic methods, embodied in one or more
software programs, that identify open reading frames (ORFs)
using a variety of heuristics, such as GRAIL, DICTION, and
GENEFINDER. Comparative sequence analysis similarly can be

25    performed using any of a variety of known programs that
identify regions with lower sequence variability.

          As further described in Example 1, below, gene
finding software programs yield a range of results. For
the newly accessioned human genomic sequence input in

30    Example 1, for example, GRAIL identified the greatest
percentage of genomic sequence as putative coding region,
2% of the data analyzed; GENEFINDER was second, calling 1%;
and DICTION yielded the least putative coding region, with
0.8% of genomic sequence called as coding region.

35          Increased reliability can be obtained when

26

consensus is required among several such methods.  Although discussed herein particularly with respect to exon calling, consensus among methods will in general increase reliability of predicting other functions as well.

5          Thus, as indicated by query 26, sequence processing 25, optionally with preprocessing 24, can be repeated with a different method, with consensus among such iterations determined and reported in process 27.

          Process 27 compares the several outputs for a
10   given input genomic sequence and identifies consensus among the separately reported results.  The consensus itself, as well as the sequence meeting that consensus, is then stored in process 29a, displayed in process 29b, and/or output to process 300 for subsequent identification of a subset
15   thereof suitable for assay.

          Multiple levels of consensus can be calculated and reported by process 27.  For example, as further described in Example 1, *infra*, process 27 can report consensus as between all specific pairs of methods of gene
20   prediction, as consensus among any one or more of the pairs of methods of gene prediction, or as among all of the gene prediction algorithms used.  Thus, in Example 1, process 27 reported that GRAIL and GENEFINDER programs agreed on 0.7% of genomic sequence, that GRAIL and DICTION agreed on 0.5%
25   of genomic sequence, and that the three programs together agreed on 0.25% of the data analyzed.  Put another way, 0.25% of the genomic sequence was identified by all three of the programs as containing putative coding region.

          Furthermore, consensus can be required among
30   different approaches to identifying a chosen function.

          For example, if the function desired to be identified is coding of protein sequence, and a first used approach to exon calling is gene prediction, the process can be repeated on the same input sequence, or subset
35   thereof, with another approach, such as comparative

sequence analysis. In such a case, where comparative
sequence analysis follows gene prediction, the comparison
can be performed not only on genomic nucleic acid sequence,
but additionally or alternatively can be performed on the
5   predicted amino acid sequence translated from the ORFs
prior identified by the gene prediction approach.

Although shown as an iterative process, the
multiple analyses required to achieve consensus can be done
in series, in parallel, or some combination thereof.

10      Predicted functional sequence, optionally
representing a consensus among a plurality of methods and
approaches for determination thereof, is passed to process
300 for identification of a subset thereof for functional
assay.

15      In the preferred embodiment of the methods of the
present invention, wherein the function sought to be
identified is protein coding, process 300 is used to
identify a subset thereof suitable for experimental
verification by physical and/or bioinformatic approaches.

20      For example, putative ORFs identified in process
200 can be classified, or binned, bioinformatically into
putative genes. This binning can be based inter alia upon
consideration of the average number of exons/gene in the
species chosen for analysis, upon density of exons that
25  have been called on the genomic sequence, and other
empirical rules. Thereafter, one or more among the gene-
specific ORFs can be chosen for subsequent use in gene
expression assay.

Where such subsequent gene expression assay uses
30  amplified nucleic acid, considerations such as desired
amplicon length, primer synthesis requirements, putative
exon length, sequence GC content, existence of possible
secondary structure, and the like can be used to identify
and select those ORFs that appear most likely successfully
35  to amplify. Where subsequent gene expression assay relies

28

upon nucleic acid hybridization, whether or not using
amplified product, further considerations involving
hybridization stringency can be applied to identify that
subset of sequences that will most readily permit sequence-
5   specific discrimination at a chosen hybridization and wash
stringency.  One particular such consideration is avoidance
of putative exons that span repetitive sequence; such
sequence can hybridize spuriously to nonspecific message,
reducing specific signal in the hybridization.
10          For bioinformatic assay, there are fewer
constraints on the sequences that can be tested
experimentally, and in this latter case therefore process
300 can output the entirety of the input sequence.
            The subset of sequences identified by process 300
15  as suitable for use in assay is then used in process 400 to
create the physical and/or informational substrate for
experimental verification of the predictions made in
process 200, and thereafter to assay those substrates.
            As mentioned, the methods of the present
20  invention are particularly useful for identifying potential
coding regions within genomic sequence.  In a preferred
embodiment of process 400, therefore, the expression of the
sequences predicted to encode protein is verified.  The
combination of the predictive and experimental methods
25  provides a powerful gene discovery engine.
            Thus, in another aspect, the present invention
provides methods and apparatus for verifying the expression
of putative genes identified within genomic sequence.  In
particular, the invention provides a novel method of
30  verifying gene expression in which expression of predicted
ORFs is measured and confirmed using a novel type of
nucleic acid microarray, the genome-derived single exon
nucleic acid microarrays of the present invention.
            Putative ORFs as predicted by a consensus of gene
35  calling, particularly gene prediction, algorithms in

process 200, and as further identified as suitable by
process 300, are amplified from genomic DNA using the
polymerase chain reaction (PCR). Although PCR is
conveniently used, other amplification approaches can also
5   be used.

Amplification schemes can be designed to capture
the entirety of each predicted ORF in an amplicon with
minimal additional (that is, intronic or intergenic)
sequence. Because ORFs predicted from human genomic
10  sequence using the methods of the present invention differ
in length, such an approach results in amplicons of varying
length.

However, most predicted ORFs are shorter than 500
bp in length, and although amplicons of at least about 100
15  or 200 base pairs can be immobilized as probes on nucleic
acid microarrays, early experimental results using the
methods of the present invention have suggested that longer
amplicons, at least about 400 or 500 base pairs, are more
effective. Furthermore, certain advantages derive from
20  application to the microarray of amplicons of defined size.

Therefore, amplification schemes can
alternatively, and preferably, be designed to amplify
regions of defined size, preferably at least about 300, 400
or 500 bp, centered about each predicted ORF. Such an
25  approach results in a population of amplicons of limited
size diversity, but that typically contain intronic and/or
intergenic nucleic acid in addition to putative ORF.

Conversely, somewhat fewer than 10% of ORFs
predicted from human genomic sequence according to the
30  methods of the present invention exceed 500 bp in length.
Portions of such extended ORFs, preferably at least about
300,400 or 500 bp in length, can be amplified. However, it
has been discovered that the percentage success at
amplifying pieces of such ORFs is low, and that such
35  putative exons are more effectively amplified when larger

30

fragments, at least about 1000 or 1500 bp, and even as large as 2000 bp are amplified.

The putative ORFs selected in process 300 are thus input into one or more primer design programs, such as

5    PRIMER3 (available online for use at http://www-genome.wi.mit.edu/cgi-bin/primer/ ), with a goal of amplifying at least about 500 base pairs of genomic sequence centered within or about ORFs predicted to be no more than about 500 bp, or at least about 1000 - 1500 bp of

10   genomic sequence for ORFs predicted to exceed 500 bp in length, and the primers synthesized by standard techniques. Primers with the requisite sequences can be purchased commercially or synthesized by standard techniques.

Conveniently, a first predetermined sequence can

15   be added commonly to the ORF-specific 5' primer and a second, typically different, predetermined sequence commonly added to each 3' ORF-unique primer.  This serves to immortalize the amplicon, that is, serves to permit further amplification of any amplicon using a single set of

20   primers complementary respectively to the common 5' and common 3' sequence elements.  The presence of these "universal" priming sequences further facilitates later sequence verification, providing a sequence common to all amplicons at which to prime sequencing reactions.  The

25   common 5' and 3' sequences further serve to add a cloning site should any of the ORFs warrant further study.

Such predetermined sequence is usefully at least about 10, 12 or 15 nt in length, and usually does not exceed about 25 nt in length.  The "universal" priming

30   sequences used in the examples presented *infra* were each 16 nt long.

The genomic DNA to be used as substrate for amplification will come from the eukaryotic species from which the genomic sequence data had originally been

35   obtained, or a closely related species, and can

31

conveniently be prepared by well known techniques from
somatic or germline tissue or cultured cells of the
organism.   See, *e.g.*, Short Protocols in Molecular Biology
: A Compendium of Methods from Current Protocols in
Molecular Biology, Ausubel et al. (eds.), 4[th] edition
(April 1999), John Wiley & Sons (ISBN: 047132938X) and
*Maniatis et al.*, Molecular Cloning : A Laboratory Manual,
2[nd] edition (December 1989), Cold Spring Harbor Laboratory
Press (ISBN: 0879693096).  Many such prepared genomic DNAs
are available commercially, with the human genomic DNAs
additionally having certification of donor informed
consent.

       Although the intronic and intergenic material
flanking putative coding regions in the amplicons could
potentially interfere with hybridizations during microarray
experiments, we have found, surprisingly, that differential
expression ratios are not significantly affected.  Rather,
the predominant effect of exon size is to alter the
absolute signal intensity, rather than its ratio.  Equally
surprising, the art had suggested that single exon probes
would not provide sufficient signal intensity for high
stringency hybridization analyses; we find that such probes
not only provide adequate signal, but have substantial
advantages, as herein described.

          After partial purification, as by size exclusion
spin column, with or without confirmation as to amplicon
quality as by gel electrophoresis, each amplicon (single
exon probe) is disposed in an array upon a support
substrate.

          Methods for creating microarrays by deposition
and fixation of nucleic acids onto support substrates are
well known in the art (Reviewed by Schena et al., see
above).

          Typically, the support substrate will be glass,
although other materials, such as amorphous or crystalline

silicon or plastics. Such plastics include polymethylacrylic, polyethylene, polypropylene, polyacrylate, polymethylmethacrylate, polyvinylchloride, polytetrafluoroethylene, polystyrene, polycarbonate,

5   polyacetal, polysulfone, celluloseacetate, cellulosenitrate, nitrocellulose, or mixtures thereof, can also be used.  Typically, the support will be rectangular, although other shapes, particularly circular disks and even spheres, present certain advantages.  Particularly

10  advantageous alternatives to glass slides as support substrates for array of nucleic acids are optical discs, as described in WO 98/12559.

The amplified nucleic acids can be attached covalently to a surface of the support substrate or, more

15  typically, applied to a derivatized surface in a chaotropic agent that facilitates denaturation and adherence by presumed noncovalent interactions, or some combination thereof.

Robotic spotting devices useful for arraying

20  nucleic acids on support substrates can be constructed using public domain specifications (The MGuide, version 2.0, http://cmgm.stanford.edu/pbrown/mguide/index.html), or can conveniently be purchased from commercial sources (MicroArray GenII Spotter and MicroArray GenIII Spotter,

25  Molecular Dynamics, Inc., Sunnyvale, CA).  Spotting can also be effected by printing methods, including those using ink jet technology.

As is well known in the art, microarrays typically also contain immobilized control nucleic acids.

30  For controls useful in providing measurements of background signal for the genome-derived single exon microarrays of the present invention, a plurality of *E. coli* genes can readily be used.  As further described in Example 1, 16 or 32 *E. coli* genes suffice to provide a robust measure of

35  background noise in such microarrays.

As is well known in the art, the amplified product disposed in arrays on a support substrate to create a nucleic acid microarray can consist entirely of natural nucleotides linked by phosphodiester bonds, or

5   alternatively can include either nonnative nucleotides, alternative internucleotide linkages, or both, so long as complementary binding can be obtained in the hybridization. If enzymatic amplification is used to produce the immobilized probes, the amplifying enzyme will impose

10  certain further constraints upon the types of nucleic acid analogs that can be generated.

Although particularly described herein as using high density microarrays constructed on planar substrates, the methods of the present invention for confirming the

15  expression of ORFs predicted from genomic sequence can use any of the known types of microarrays, as herein defined, including lower density planar arrays, and microarrays on nonplanar, nonunitary, distributed substrates.

For example, gene expression can be confirmed

20  using hybridization to lower density arrays, such as those constructed on membranes, such as nitrocellulose, nylon, and positively-charged derivatized nylon membranes. Further, gene expression can also be confirmed using nonplanar, bead-based microarrays such as are described in

25  Brenner *et al.*, *Proc. Natl. Acad. Sci. USA* 97(4):166501670 (2000); U.S. Patent No. 6,057,107; and U.S. Patent No. 5,736,330. In theory, a packed collection of such beads provides in aggregate a higher density of nucleic acid probe than can be achieved with spotting or lithography

30  techniques on a single planar substrate.

Planar microarrays on solid substrates, however, provide certain useful advantages, including high throughput and compatibility with existing readers. For example, each standard microscope slide can include at

35  least 1000, typically at least 2000, preferably 5000 and

upto 10,000 – 50,000 or more nucleic acid probes of discrete sequence. The number of sequences deposited will depend on their required application.

Each putative gene can be represented in the array by a single predicted ORF. Alternatively, genes can be represented by more than one predicted ORF. For purposes of measuring differential splicing, more than one predicted ORF will be provided for a putative gene. And as is well known in the art, each probe of defined sequence, representing a single predicted ORF, can be deposited in a plurality of locations on a single microarray to provide redundancy of signal.

The genome-derived single exon microarrays described above differ in several fundamental and advantageous ways from microarrays presently used in the gene expression art, including (1) those created by deposition of mRNA-derived nucleic acids, (2) those created by *in situ* synthesis of oligonucleotide probes, and (3) those constructed from yeast genomic DNA.

Most nucleic acid microarrays that are in use for study of eukaryotic gene expression have as immobilized probes nucleic acids that are derived – either directly or indirectly – from expressed message. As discussed above, it is common, for example, for such microarrays to be derived from cDNA/EST libraries, either from those previously described in the literature, see Lennon *et al.*, or from the *de novo* construction of "problem specific" libraries targeted at a particular biological question, R.S. Thomas *et al.*, *Cancer Res.* (in press). Such microarrays are herein collectively denominated "EST microarrays".

Such EST microarrays by definition can measure expression only of those genes found in EST libraries, shown herein to represent only a fraction of expressed genes. Furthermore, such libraries – and thus microarrays

35

based thereupon — are biased by the tissue or cell type of
message origin, by the expression levels of the respective
genes within the tissues, and by the ability of the message
successfully to have been reverse-transcribed and cloned.

5          Thus, as further discussed in Example 1, the
methods of the present invention enable sequences that do
not appear in EST or other expression databases to be
determined - subsequently arrayed for expression
measurements could not, therefore, have been represented as
10  probes on an EST microarray.  And as further demonstrated
in the examples, *infra*, the remaining population of genes
identified from genomic sequence by the methods of the
present invention — that is, the one third of sequences
that had previously been accessioned in EST or other
15  expression databases — are biased toward genes with higher
expression levels.

          Representation of a message in an EST and/or cDNA
library depends upon the successful reverse transcription,
optionally but typically with subsequent successful
20  cloning, of the message.  This introduces substantial bias
into the population of probes available for arraying in EST
microarrays.

          In contrast, neither reverse transcription nor
cloning is required to produce the probes arrayed on the
25  genome-derived single exon microarrays of the present
invention.  And although the ultimate deposition of a probe
on the genome-derived single exon microarray of the present
invention depends upon a successful amplification from
genomic material, *a priori* knowledge of the sequence of the
30  desired amplicon affords greater opportunity to recover any
given probe sequence recalcitrant to amplification than is
afforded by the requirement for successful reverse
transcription and cloning of unknown message in EST
approaches.

35          Thus, the genome-derived single exon microarrays

36

of the present invention present a far greater diversity of probes for measuring gene expression, with far less bias, than do EST microarrays presently used in the art.

5          As a further consequence of their ultimate origin from expressed message, the probes in EST microarrays often contain poly-A (or complementary poly-T) stretches derived from the poly-A tail of mature mRNA. These homopolymeric stretches contribute to cross-hybridization, that is, to a spurious signal occasioned by hybridization to the

10     homopolymeric tail of a labeled cDNA that lacks sequence homology to the gene-specific portion of the probe.

          In contrast, the probes arrayed in the genome-derived single exon microarrays of the present invention lack homopolymeric stretches derived from message

15     polyadenylation, and thus can provide more specific signal. Typically, at least about 50, 60 or 75% of the probes on the genome-derived single exon microarrays of the present invention lack homopolymeric regions consisting of A or T, where a homopolymeric region is defined for purposes herein

20     as stretches of 25 or more, typically 30 or more, identical nucleotides.

          A further distinction, which also affects the specificity of hybridization, is occasioned by the typical derivation of EST microarray probes from cloned material.

25     Because much of the probe material disposed as probes on EST microarrays is excised or amplified from plasmid, phage, or phagemid vectors, EST microarrays typically include a fair amount of vector sequence, more so when the probes are amplified, rather than excised, from the vector.

30          In contrast, the vast majority of probes in the genome-derived single exon microarrays of the present invention contain no prokaryotic or bacteriophage vector sequence, having been amplified directly or indirectly from genomic DNA. Typically, therefore, at least about 50, 60,

35     70 or 80% or more of individual exon-including probes

37

disposed on a genome-derived single exon microarray of the
present invention lack vector sequence, and particularly
lack sequences drawn from plasmids and bacteriophage.
Preferably, at least about 85, 90 or more than 90% of exon-
5    including probes in the genome-derived single exon
microarray of the present invention lack vector sequence.
With attention to removal of vector sequences through
preprocessing 24, percentages of vector-free exon-including
probes can be as high as 95 - 99%.  The substantial absence
10   of vector sequence from the genome-derived single exon
microarrays of the present invention results in greater
specificity during hybridization, since spurious cross-
hybridization to a probe vector sequence is reduced.

          As a further consequence of excision or
15   amplification of probes from vectors in construction of EST
microarrays, the probes arrayed thereon often contain
artificial sequence, derived from vector polylinker
multiple cloning sites, at both 5' and 3' ends.  The probes
disposed upon the genome-derived single exon microarrays
20   need have no such artificial sequence appended thereto.

          As mentioned above, however, the ORF-specific
primers used to amplify putative ORFs can include
artificial sequences, typically 5' to the ORF-specific
primer sequence, useful for "universal" (that is,
25   independent of ORF sequence) priming of subsequent
amplification or sequencing reactions.  When such
"universal" 5' and/or 3' priming sequences are appended to
the amplification primers, the probes disposed upon the
genome-derived single exon microarray will include
30   artificial sequence similar to that found in EST
microarrays.  However, the genome-derived single exon
microarray of the present invention can be made without
such sequences, and if so constructed, presents an even
smaller amount of nonspecific sequence that would
35   contribute to nonspecific hybridization.

38

Yet another consequence of typical use of cloned
material as probes in EST microarrays is that such
microarrays contain probes that result from cloning
artifacts, such as chimeric molecules containing coding

5    region of two separate genes.  Derived from genomic
material, typically not thereafter cloned, the probes of
the genome-derived single exon microarrays of the present
invention lack such cloning artifacts, and thus provide
greater specificity of signal in gene expression

10   measurements.

A further consequence of the cloned origin of
probes on many EST microarrays is that the individual
probes often have disparate sizes, which can cause the
optimal hybridization stringency to vary among probes on a

15   single microarray.  In contrast, as discussed above, the
probes arrayed on the genome-derived single exon
microarrays of the present invention can readily be
designed to have a narrow distribution in sizes, with the
range of probe sizes no greater than about 10% of the

20   average size, typically no greater than about 5% of the
average probe size.

Because of their origin from fully- or partially-
spliced message, probes disposed upon EST arrays will often
include multiple exons.  The percentage of such exon-

25   spanning probes in an EST microarray can be calculated, on
average, based upon the predicted number of exons/gene for
the given species and the average length of the immobilized
probes.  For human genes, the near-complete sequence of
human chromosome 22, Dunham *et al.*, *Nature* 402(6761):489-95

30   (1999), predicts that human genes average 5.5 exons/gene.
Even with probes of 200 - 500 bp, the vast majority of
human EST microarray probes include more than one exon.

In contrast, by virtue of their origin from
algorithmically identified ORFs in genomic sequence, the

35   probes in the genome-derived single exon microarrays of the

present invention can consist of individual exons.  Thus,
in contrast to EST microarrays, at least about 50, 60, 70,
75, 80, 85, 95 or 99% of probes deposited in the genome-
derived microarray of the present invention consist of, or
5    include, no more than one predicted ORF.

This provides the ability, not readily achieved
using EST microarrays, to use the genome-derived single
exon microarrays of the present invention to measure
tissue-specific expression of individual exons, which in
10   turn allows differential splicing events to be detected and
characterized, and in particular, allows the correlation of
differential splicing to tissue-specific expression
patterns.

Furthermore, the exons that are represented in
15   EST microarrays are often biased toward the 3' or 5' end of
their respective genes, since sequencing strategies used
for EST identification are so biased.  In contrast, no such
3' or 5' bias necessarily inheres in the selection of exons
for disposition on the genome-derived single exon
20   microarrays of the present invention.

Conversely, the probes provided on the genome-
derived single exon microarrays of the present invention
typically, but need not necessarily, include intronic
and/or intergenic sequence that is absent from EST
25   microarrays, which are derived from mature mRNA.
Typically, at least about 50, 60, 70, 80 or 90% of the
exon-including probes on the genome-derived single exon
microarrays of the present invention include sequence drawn
from noncoding regions. As discussed above, the additional
30   presence of noncoding region does not significantly
interfere with measurement of gene expression, and provides
the additional opportunity to assay prespliced RNA, and
thus measure such phenomena such as nuclear export control.

The genome-derived single exon microarrays of the
35   present invention are also quite different from *in situ*

40

synthesis microarrays, where probe size is severely
constrained by inadequacies in the photolithographic
synthesis process.

 5      Typically, probes arrayed on *in situ* synthesis
microarrays are limited to a maximum of about 25 bp.  As a
well known consequence, hybridization to such chips must be
performed at low stringency.  In order, therefore, to
achieve unambiguous sequence-specific hybridization
results, the *in situ* synthesis microarray requires
10   substantial redundancy, with concomitant programmed
arraying for each probe of probe analogues with altered
(*i.e.*, mismatched) sequence.

        In contrast, the longer probe length of the
genome-derived single exon microarrays of the present
15   invention allows much higher stringency hybridization and
wash.  Typically, therefore, exon-including probes on the
genome-derived single exon microarrays of the present
invention average at least about 100, 200, 300, 400 or
500 bp in length.  By obviating the need for substantial
20   probe redundancy, this approach permits a higher density of
probes for discrete exons or genes to be arrayed on the
microarrays of the present invention than can be achieved
for *in situ* synthesis microarrays.

        A further distinction is that the probes in *in*
25   *situ* synthesis microarrays typically are covalently linked
to the substrate surface.  In contrast, the probes disposed
on the genome-derived microarray of the present invention
typically are, but need not necessarily be, bound
noncovalently to the substrate.

30      Furthermore, the short probe size on *in situ*
microarrays causes large percentage differences in the
melting temperature of probes hybridized to their
complementary target sequence, and thus causes large
percentage differences in the theoretically optimum
35   stringency across the array as a whole.

In contrast, the larger probe size in the microarrays of the present invention create lower percentage differences in melting temperature across the range of arrayed probes.

5      A further significant advantage of the microarrays of the present invention over *in situ* synthesized arrays is that the quality of each individual probe can be confirmed before deposition. In contrast, the quality of probes cannot be assessed on a probe-by-probe

10     basis for the *in situ* synthesized microarrays presently being used.

The genome-derived single exon microarrays of the present invention are also distinguished over, and present substantial benefits over, the genome-derived microarrays

15     from lower eukaryotes such as yeast. Lashkari *et al.*, *Proc. Natl. Acad. Sci. USA* 94:13057-13062 (1997).

Only about 220 - 250 of the 6100 or so nuclear genes in *Saccharomyces cerevisiae* — that is, only about 4 - 5% — have standard, spliceosomal, introns, Lopez *et al.*,

20     *Nucl. Acids Res.* 28:85-86 (2000); Spingola *et al.*, *RNA* 5(2):221-34 (1999). Furthermore, the entire yeast genome has already been sequenced. These two facts permit the ready amplification and disposition of single-ORF amplicons on such microarray without the requirement for antecedent

25     use of gene prediction and/or comparative sequence analyses.

Thus, a significant aspect of the present invention is the ability to identify and to confirm expression of predicted coding regions in genomic sequence

30     drawn from eukaryotic organisms that have a higher percentage of genes having introns than do yeast such as *Saccharomyces cerevisiae*, particularly in genomic sequence drawn from eukaryotes in which at least about 10, 20 or 50% of protein-encoding genes have introns. In preferred

35     embodiments, the methods and apparatus of the present

42

invention are used to identify and confirm expression of
novel genes from genomic sequence of eukaryotes in which
the average number of introns per gene is at least about
one, two or three or more.

5          After the physical substrate is prepared,
experimental verification of predicted function is
performed.

           In a preferred embodiment of the present
invention, where the function sought to be identified in
10    genomic sequence is protein coding, experimental
verification is performed by measuring expression of the
putative ORFs, typically through nucleic acid hybridization
experiments, and in particularly preferred embodiments,
through hybridization to genome-derived single exon
15    microarrays prepared as above- described.

           Expression is conveniently measured and expressed
for each probe in the microarray as a ratio of the
expression measured concurrently in a plurality of mRNA
sources, according to techniques well known in the
20    microarray art, *Reviewed in* Schena et al., and as further
described in Example 2, below.  The mRNA source for the
reference against which specific expression is measured can
be drawn from a homogeneous mRNA source, such as a single
cultured cell-type, or alternatively can be heterogeneous,
25    as from a pool of mRNA derived from multiple tissues and/or
cell types, as further described in Example 2, *infra*.

           mRNA can be prepared by standard techniques, see
Ausubel et al. and Maniatis *et al.*, or purchased
commercially.  The mRNA is then typically reverse-
30    transcribed in the presence of labeled nucleotides: the
index source (that in which expression is desired to be
measured) is reverse transcribed in the presence of
nucleotides labeled with a first label, typically a
fluorophore (fluorochrome; fluor; fluorescent dye); the
35    reference source is reverse transcribed in the presence of

43

a second label, typically a fluorophore, typically
fluorometrically-distinguishable from the first label.  As
further described in Example 2, *infra*, Cy3 and Cy5 dyes
prove particularly useful in these methods.  After partial
5  purification of the index and reference targets,   .
hybridization to the probe array is conducted according to
standard techniques, typically under a coverslip.

    After wash, microarrays are conveniently scanned
using a commercial microarray scanning device, such as a
10  Gen3 Scanner (Molecular Dynamics, Sunnyvale, CA).  Data on
expression is then passed, with or without interim storage,
to process 500, where the results for each probe are
related to the original sequence.

    Often, hybridization of target material to the
15  genome-derived single exon microarray will identify certain
of the probes thereon as of particular interest.  Thus, it
is often desirable that the user be able readily to obtain
sufficient quantities of an individual probe, either for
subsequent arrayed deposition upon an additional support
20  substrate, often as part of a microarray having a plurality
of probes so identified, or alternatively or additionally
as a solitary solid-phase or solution-phase probe, for
further use.

    Thus, in another aspect, the present invention
25  provides compositions and kits for the ready production of
nucleic acids identical in sequence to, or substantially
identical in sequence to, probes on the genome-derived
single exon microarrays of the present invention.

    In this aspect, a small quantity of each probe is
30  disposed, typically without attachment to substrate, in a
spatially-addressable ordered set, typically one per well
of a microtiter dish.  Although a 96 well microtiter plate
can be used, greater efficiency is obtained using higher
density arrays, such as are provided by microtiter plates
35  having 384, 864, 1536, 3456, 6144, or 9600 wells, and

44

although microtiter plates having physical depressions
(wells) are conveniently used, any device that permits
addressable withdrawal of reagent from fluidly-
noncommunicating areas can be used.

5          In this aspect of the invention, therefore, a
fluidly noncommunicating addressable ordered set of
individual probes, corresponding to those on a genome-
derived single exon microarray, is provided, with each
probe in sufficient quantity to permit amplification, such
10   as by PCR.  As earlier mentioned, the ORF-specific
5' primers used for genomic amplification can have a first
common sequence added thereto, and the ORF-specific 3'
primers used for genomic amplification can have a second,
different, common sequence added thereto, thus permitting,
15   in this preferred embodiment, the use of a single set of 5'
and 3' primers to amplify any one of the probes from the
amplifiable ordered set.

Each discrete amplifiable probe can also be
packaged with amplification primers, solutes, buffers,
20   etc., and can be provided in dry (e.g., lyophilized) form
or wet, in the latter case typically with addition of
agents that retard evaporation.

In another aspect of the present invention, a
genome-derived single-exon microarray is packaged together
25   with such an ordered set of amplifiable probes
corresponding to the probes, or one or more subsets of
probes, thereon.  In alternative embodiments, the ordered
set of amplifiable probes is packaged separately from the
genome-derived single exon microarray.

30          In some embodiments, the microarray and/or
ordered probe set are further packaged with recordable
media that provide probe identification and addressing
information, and that can additionally contain annotation
information, such as gene expression data.  Such recordable
35   media can be packaged with the microarray, with the ordered

45

probe set, or with both.

If the microarray is constructed on a substrate that incorporates recordable media, such as is described in international patent application no. WO 98/12559, then
5   separate packaging of the genome-derived single exon microarray and the bioinformatic information is not required.

The amount of amplifiable probe material should be sufficient to permit at least one amplification
10  sufficient for subsequent hybridization assay.

Although the use of high density genome-derived microarrays on solid planar substrates is presently a preferred approach for the physical confirmation and characterization of the expression of sequences predicted
15  to encode protein, other types of microarrays (as herein defined) can also be used.

Furthermore, as earlier mentioned, experimental verification of the function predicted from genomic sequence in process 200 can be bioinformatic, rather than,
20  or additional to, physical verification.

For example, where the function desired to be identified is protein coding, the predicted ORFs can be compared bioinformatically to sequences known or suspected of being expressed.
25          Thus, the sequences output from process 300 (or process 200), can be used to query expression databases, such as EST databases, SNP ("single nucleotide polymorphism") databases, known cDNA and mRNA sequences, SAGE ("serial analysis of gene expression") databases, and
30  more generalized sequence databases that allow query for expressed sequences. Such query can be done by any sequence query algorithm, such as BLAST ("basic local alignment search tool"). The results of such query — including information on identical sequences and
35  information on nonidentical sequences that have diffuse or

46

focal regions of sequence homology to the query sequence —
can then be passed directly to process 500, or used to
inform analyses subsequently undertaken in process 200,
process 300, or process 400.

5          Experimental data, whether obtained by physical
or bioinformatic assay in process 400, is passed to process
500 where it is usefully related to the sequence data
itself, a process colloquially termed "annotation". Such
annotation can be done using any technique that usefully
10   relates the functional information to the sequence, as, for
example, by incorporating the functional data into the
record itself, by linking records in a hierarchical or
relational database, by linking to external databases, or
by a combination thereof. Such database techniques are
15   well within the skill in the art.

The annotated sequence data can be stored
locally, uploaded to genomic sequence database 100, and/or
displayed 800.

The methods and apparatus of the present
20   invention rapidly produce functional information from
genomic sequence. Coupled with the escalating pace at
which sequence now accumulates, the rapid pace of sequence
annotation produces a need for methods of displaying the
information in meaningful ways.

25          FIG. 3 shows visual display 80 presenting a
single genomic sequence annotated according to the present
invention. Because of its nominal resemblance to artistic
works of Piet Mondrian, visual display 80 is alternatively
described herein as a "Mondrian".

30          Each of the visual elements of display 80 is
aligned with respect to the genomic sequence being
annotated (hereinafter, the "annotated sequence"). Given
the number of nucleotides typically represented in an
annotated sequence, representation of individual
35   nucleotides would rarely be readable in hard copy output of

47

display 80. Typically, therefore, the annotated sequence
is schematized as rectangle 89, extending from the left
border of display 80 to its right border. By convention
herein, the left border of rectangle 89 represents the

5   first nucleotide of the sequence and the right border of
rectangle 89 represents the last nucleotide of the
sequence.

        As further discussed below, however, the Mondrian
visual display of annotated sequence can serve as a

10  convenient graphical user interface for computerized
representation, analysis, and query of information stored
electronically. For such use, the individual nucleotides
can conveniently be linked to the X axis coordinate of
rectangle 89. This permits the annotated sequence at any

15  point within rectangle 89 readily to be viewed, either
automatically — for example, by time-delayed appearance of
a small overlaid window upon movement of a cursor or other
pointer over rectangle 89 — or through user intervention,
as by clicking a mouse or other pointing device at a point

20  in rectangle 89.

        Visual display 80 is generated after user
specification of the genomic sequence to be displayed.
Such specification can consist of or include an accession
number for a single clone (e.g., a single BAC accessioned

25  into GenBank), wherein the starting and stopping
nucleotides are thus absolutely identified, or
alternatively can consist of or include an anchor or
fulcrum point about which a chosen range of sequence is
anchored, thus providing relative endpoints for the

30  sequence to be displayed. For example, the user can anchor
such a range about a given chromosomal map location, gene
name, or even a sequence returned by query for similarity
or identity to an input query sequence. When visual
display 80 is used as a graphical user interface to

35  computerized data, additional control over the first and

48

last displayed nucleotide will typically be dynamically
selectable, as by use of standard zooming and/or selection
tools.

        Field 81 of visual display 80 is used to present
5   the output from process 200, that is, to present the
bioinformatic prediction of those sequences having the
desired function within the genomic sequence.  Functional
sequences are typically indicated by at least one rectangle
83 (83a, 83b, 83c), the left and right borders of which
10  respectively indicate, by their X-axis coordinates, the
starting and ending nucleotides of the region predicted to
have function.

        Where a single bioinformatic method or approach
identifies a plurality of regions having the desired
15  function, a plurality of rectangles 83 is disposed
horizontally in field 81.  Where multiple methods and/or
approaches are used to identify function, each such method
and/or approach can be represented by its own series of
horizontally disposed rectangles 83, each such horizontally
20  disposed series of rectangles offset vertically from those
representing the results of the other methods and
approaches.

        Thus, rectangles 83a in FIG. 3 represent the
functional predictions of a first method of a first
25  approach for predicting function, rectangles 83b represent
the functional predictions of a second method and/or second
approach for predicting that function, and rectangles 83c
represent the predictions of a third method and/or
approach.

30      Where the function desired to be identified is
protein coding, field 81 is used to present the
bioinformatic prediction of sequences encoding protein.
For example, rectangles 83a can represent the results from
GRAIL or GRAIL II, rectangles 83b can represent the results
35  from GENEFINDER, and rectangles 83c can represent the

49